

# Cumulative Restricted Boltzmann Machines for Ordinal Matrix Data Analysis

Truyen Tran<sup>†‡</sup>, Dinh Phung<sup>†</sup>, Svetha Venkatesh<sup>†</sup>

<sup>†</sup>Pattern Recognition and Data Analytics, Deakin University, Australia

<sup>‡</sup>Department of Computing, Curtin University, Australia  
 {truyen.tran,dinh.phung,svetha.venkatesh}@deakin.edu.au

## Abstract

Ordinal data is omnipresent in almost all multiuser-generated feedback - questionnaires, preferences etc. This paper investigates modelling of ordinal data with Gaussian restricted Boltzmann machines (RBMs). In particular, we present the model architecture, learning and inference procedures for both vector-variate and matrix-variate ordinal data. We show that our model is able to capture latent opinion profile of citizens around the world, and is competitive against state-of-art collaborative filtering techniques on large-scale public datasets. The model thus has the potential to extend application of RBMs to diverse domains such as recommendation systems, product reviews and expert assessments.

## 1 Introduction

Restricted Boltzmann machines (RBMs) [36, 9, 20] have recently attracted significant interest due to their versatility in a variety of unsupervised and supervised learning tasks [35, 18, 25], and in building deep architectures [14, 31]. A RBM is a bipartite undirected model that captures the generative process in which a data vector is generated from a binary hidden vector. The bipartite architecture enables very fast data encoding and sampling-based inference; and together with recent advances in learning procedures, we can now process massive data with large models [13, 37, 2].

This paper presents our contributions in developing RBM specifications as well as learning and inference procedures for multivariate ordinal data. This extends and consolidates the reach of RBMs to a wide range of user-generated domains - social responses, recommender systems, product/paper reviews, and expert assessments of health and ecosystems indicators. Ordinal variables are qualitative in nature – the absolute numerical assignments are not important but the relative order is. This renders numerical transforms and real-valued treatments inadequate. Current RBM-based treatments, on the other hand, ignore the ordinal nature and treat data as unordered categories [35, 40]. While convenient, this has several drawbacks: First, order information is not utilised, leading to more parameters than necessary - each category needs parameters. Second, since categories are considered independently, it is less interpretable in terms of how ordinal levels are generated. Better modelling should account for the ordinal generation process.

Adapting the classic idea from [24], we assume that each ordinal variable is generated by an underlying latent utility, along with a threshold per ordinal level. As soon as the utility passes

the threshold, its corresponding level is selected. As a result, this process would implicitly encode the order. Our main contribution here is a novel RBM architecture that accounts for multivariate, ordinal data. More specifically, we further assume that the latent utilities are Gaussian variables connected to a set of binary hidden factors (i.e., together they form a Gaussian RBM [14]). This offers many advantages over the standard approach that imposes a fully connected Gaussian random field over utilities [17, 15]: First, utilities are seen as being generated from a set of binary factors, which in many cases represent the user’s hidden profile. Second, utilities are decoupled given the hidden factors, making parallel sampling easier. And third, the posteriors of binary factors can be estimated from the ordinal observations, facilitating dimensionality reduction and visualisation. We term our model Cumulative RBM (CRBM)<sup>1</sup>.

This new model behaves differently from standard Gaussian RBMs since utilities are never observed in full. Rather, when an ordinal level of an input variable is observed, it poses an *interval constraint* over the corresponding utility. The distribution over the utilities now becomes a *truncated* multivariate Gaussian. This also has another consequence during learning: While in standard RBMs we need to sample for the *free-phase* only (e.g., see [13]), now we also need to sample for the *clamped-phase*. As a result, we introduce a double persistent contrastive divergence (PCD) learning procedure, as opposed to the single PCD in [37].

The second contribution is in advancing these ordinal RBMs from modelling i.i.d. vectors to modelling matrices of correlated entries. These ordinal matrices are popular in multiuser-generated assessments: Each user would typically judge a number of items producing a user-specific data vector where *intra-vector* entries are inherently correlated. Since user’s choices are influenced by their peers, these *inter-vector* entries are no longer independent. The idea is borrowed from a recent work in [40] which models both the user-specific and item-specific processes. More specifically, an ordinal entry is assumed to be jointly generated from user-specific latent factors and item-specific latent factors. This departs significantly from the standard RBM architecture: we no longer map from a visible vector to an hidden vector but rather map from a visible matrix to two hidden matrices.

In experiments, we demonstrate that our proposed CRBM is capable of capturing the latent profile of citizens around the world. Our model is also competitive against state-of-the-art collaborative filtering methods on large-scale public datasets.

We start with the RBM structure for ordinal vectors in Section 2, and end with the general structure for ordinal matrices in Section 3. Section 4 presents experiments validating our ordinal RBMs in modelling citizen’s opinions worldwide and in collaborative filtering. Section 5 discusses related work, which is then followed by the conclusions.

## 2 Cumulative RBM for Vectorial Data

### 2.1 Model Definition

Denote by  $\mathbf{v} = (v_1, v_2, \dots, v_N)$  the set of ordinal observations. For ease of presentation we assume for the moment that observations are homogeneous, i.e., observations are drawn from the same discrete ordered category set  $S = \{c_1 \prec c_2 \prec \dots \prec c_L\}$  where  $\prec$  denotes the order in some sense. We further assume that each ordinal  $v_i$  is solely generated from an underlying latent

---

<sup>1</sup>The term ‘cumulative’ is to be consistent with the statistical literature when referring to the ordinal treatment in [24].

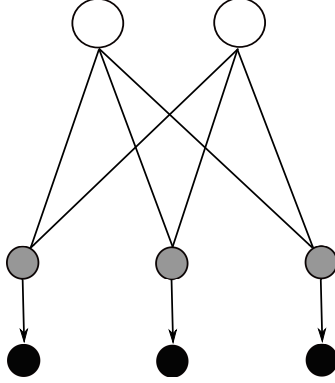


Figure 1: Model architecture of the Cumulative Restricted Boltzmann Machine (CRBM). Filled nodes represent observed ordinal variables, shaded nodes are Gaussian utilities, and empty nodes represent binary hidden factors.

utility  $u_i \in \mathbb{R}$  as follows [24]

$$P(v_i = c_l \mid u_i) = \begin{cases} \mathbb{I}[-\infty < u_i \leq \theta_{i1}] & l = 1 \\ \mathbb{I}[\theta_{i(l-1)} < u_i \leq \theta_{il}] & 1 < l \leq L-1 \\ \mathbb{I}[\theta_{i(L-1)} < u_i < \infty] & l = L \end{cases} \quad (1)$$

where  $\theta_{i1} < \theta_{i2} < \dots < \theta_{i(L-1)}$  are threshold parameters. In words, we choose an ordered category on the basis of the interval to which the underlying utility belongs.

The utilities are connected with a set of hidden *binary factors*  $\mathbf{h} = (h_1, h_2, \dots, h_K) \in \{0, 1\}^K$  so that the two layers of  $(\mathbf{u}, \mathbf{h})$  form a unidirectional bipartite graph of Restricted Boltzmann Machines (RBMs) [36, 9, 14]. Binary factors can be considered as the hidden features that govern the generation of the observed ordinal data. Thus the generative story is: we start from the binary factors to generate utilities, which, in turn, generate ordinal observations. See, for example, Fig. 1 for a graphical representation of the model.

Let  $\Psi(\mathbf{u}, \mathbf{h}) \geq 0$  be the model potential function, which can be factorised as a result of the bipartite structure as follows

$$\Psi(\mathbf{u}, \mathbf{h}) = \left[ \prod_i \phi_i(u_i) \right] \left[ \prod_{i,k} \psi_{ik}(u_i, h_k) \right] \left[ \prod_k \phi_k(h_k) \right]$$

where  $\phi_i, \psi_{ik}$  and  $\phi_k$  are local potential functions. The model joint distribution is defined as

$$P(\mathbf{v}, \mathbf{u}, \mathbf{h}) = \frac{1}{Z} \Psi(\mathbf{u}, \mathbf{h}) \prod_i P(v_i \mid u_i) \quad (2)$$

where  $Z = \int_{\mathbf{u}} \sum_{\mathbf{h}} \Psi(\mathbf{u}, \mathbf{h}) d\mathbf{u}$  is the normalising constant.

We assume the utility layer and the binary factor layer form a Gaussian RBM<sup>2</sup> [14]. This translates into the local potential functions as follows

<sup>2</sup>This is for convenience only. In fact, we can replace Gaussian by any continuous distribution in the exponential family.

$$\phi_i(u_i) = \exp \left\{ -\frac{u_i^2}{2\sigma_i^2} + \alpha_i u_i \right\}; \quad \psi_{ik}(u_i, h_k) = \exp \{w_{ik} u_i h_k\}; \quad \phi_k(h_k) = \exp \{\gamma_k h_k\} \quad (3)$$

where  $\sigma_i$  is the standard deviation of the  $i$ -th utility,  $\{\alpha_i, \gamma_k, w_{ik}\}$  are free parameters for  $i = 1, 2, \dots, N$  and  $k = 1, 2, \dots, K$ .

The ordinal assumption in Eq. (1) introduces *hard constraints* that we do not see in standard Gaussian RBMs. Whenever an ordered category  $v_i$  is observed, the corresponding utility is automatically *truncated*, i.e.,  $u_i \in \Omega(v_i)$ , where  $\Omega(v_i)$  is the new domain of  $u_i$  defined by  $v_i$  as in Eq. (1). In particular, the utility is truncated from above if the ordinal level is the lowest, from below if the level is the largest, and from both sides otherwise. For example, the conditional distribution of the latent utility  $P(u_i | v_i, \mathbf{h})$  is a truncated Gaussian

$$P(u_i | v_i, \mathbf{h}) \propto \mathbb{I}[u_i \in \Omega(v_i)] \mathcal{N}(u_i; \mu_i(\mathbf{h}), \sigma_i) \quad (4)$$

where  $\mathcal{N}(u_i; \mu_i(\mathbf{h}), \sigma_i)$  is the normal density distribution of mean  $\mu_i(\mathbf{h})$  and standard deviation  $\sigma_i$ . The mean  $\mu_i(\mathbf{h})$  is computed as

$$\mu_i(\mathbf{h}) = \sigma_i^2 \left( \alpha_i + \sum_{k=1}^K w_{ik} h_k \right) \quad (5)$$

As a generative model, we can estimate the probability that an ordinal level is being generated from hidden factors  $\mathbf{h}$  as follows

$$P(v_i = c_l | \mathbf{h}) = \int_{u_i \in \Omega(c_l)} P(u_i | \mathbf{h}) = \begin{cases} \Phi(\theta_1^*) & l = 1 \\ \Phi(\theta_l^*) - \Phi(\theta_{(l-1)}^*) & 1 < l \leq L-1 \\ 1 - \Phi(\theta_{L-1}^*) & l = L \end{cases} \quad (6)$$

where  $\theta_l^* = \frac{\theta_l - \mu_i(\mathbf{h})}{\sigma_i}$ , and  $\Phi(\cdot)$  is the cumulative distribution function of the Gaussian. Given this property, we term our model by Cumulative Restricted Boltzmann Machine (CRBM).

Finally, the thresholds are parameterised so that the lowest threshold is fixed to a constant  $\theta_{i1} = \tau_{i1}$  and the higher thresholds are spaced as  $\theta_{il} = \theta_{i(l-1)} + e^{\tau_{il}}$  with free parameter  $\tau_{il}$  for  $l = 2, 3, \dots, L-1$ .

## 2.2 Factor Posteriors

Often we are interested in the posterior of factors  $\{P(h_k | \mathbf{v})\}_{k=1}^K$  as it can be considered as a summary of the data  $\mathbf{v}$ . The nice thing is that it is now numerical and can be used for other tasks such as clustering, visualisation and prediction.

Like standard RBMs, the factor posteriors given the utilities are conditionally independent and assume the form of logistic units

$$P(h_k = 1 | \mathbf{u}) = \frac{1}{1 + \exp(-\gamma_k - \sum_i w_{ik} u_i)} \quad (7)$$

However, since the utilities are themselves hidden, the posteriors given only the ordinal observations are not independent:

$$P(h_k | \mathbf{v}) = \sum_{\mathbf{h}_{-k}} \int_{\mathbf{u} \in \Omega(\mathbf{v})} P(\mathbf{h}, \mathbf{u} | \mathbf{v}) d\mathbf{u} \quad (8)$$

where  $\mathbf{h}_{-k} = \mathbf{h} \setminus h_k$  and  $\Omega(\mathbf{v}) = \Omega(v_1) \times \Omega(v_2) \times \dots \Omega(v_N)$  is the domain of the utility constrained by  $\mathbf{v}$  (see Eq. (1)). Here we describe two approximation methods, namely Markov chain Monte Carlo (MCMC) and variational method (mean-field).

**MCMC.** We can exploit the bipartite structure of the RBM to run layer-wise Gibbs sampling: sample the truncated utilities in parallel using Eq. (4) and the binary factors using Eq. (7). Finally, the posteriors are estimated as  $P(h_k | \mathbf{v}) \approx \frac{1}{n} \sum_{s=1}^n h_k^{(s)}$  for  $n$  samples.

**Variational method.** We make the approximation

$$P(\mathbf{h}, \mathbf{u} | \mathbf{v}) \approx \prod_k Q_k(h_k) \prod_i Q_i(u_i)$$

Minimising the Kullback-Leibler divergence between  $P(\mathbf{h}, \mathbf{u} | \mathbf{v})$  and its approximation leads the following recursive update

$$Q_k(h_k^{(t+1)} = 1) \leftarrow \frac{1}{1 + \exp(-\gamma_k - \sum_i w_{ik} \langle u_i \rangle_{Q_i^{(t)}})} \quad (9)$$

$$Q_i(u_i^{(t+1)}) \leftarrow \frac{1}{\kappa_i^{(t)}} \mathbb{I}[u_i \in \Omega(v_i)] \mathcal{N}(u_i; \hat{\mu}_i(\mathbf{h}^{(t)}), \sigma_i) \quad (10)$$

where  $t$  is the update index of the recursion,  $\langle u_i \rangle_{Q_i^{(t)}}$  is the mean of utility  $u_i$  with respect to  $Q_i(u_i^{(t)})$ ,  $\kappa_i^{(t)} = \int_{u_i \in \Omega(v_i)} \mathcal{N}(u_i; \mu_i(\mathbf{h}^{(t)}), \sigma_i)$  is the normalising constant, and  $\hat{\mu}_i(\mathbf{h}^{(t)}) = \sigma_i^2 \left( \alpha_i + \sum_{k=1}^K w_{ik} Q_k(h_k^{(t)} = 1) \right)$ . Finally, we obtain  $P(h_k | \mathbf{v}) \approx Q_k(h_k = 1)$ .

## 2.3 Prediction

An important task is *prediction* of the ordinal level of an unseen variable given the other seen variables, where we need to estimate the following predictive distribution

$$P(v_j | \mathbf{v}) = \sum_{\mathbf{h}} \int_{u_j \in \Omega(v_j)} \int_{\mathbf{u} \in \Omega(\mathbf{v})} P(\mathbf{h}, u_j, \mathbf{u} | \mathbf{v}) d\mathbf{u} du_j \quad (11)$$

Unfortunately, now  $(h_1, h_2, \dots, h_K)$  are coupled due to the integration over  $\{u_j, \mathbf{u}\}$  making the evaluation intractable, and thus approximation is needed.

For simplicity, we assume that the seen data  $\mathbf{v}$  is informative enough so that  $P(\mathbf{h} | v_j, \mathbf{v}) \approx P(\mathbf{h} | \mathbf{v})$ . Thus we can rewrite Eq. (11) as

$$P(v_j | \mathbf{v}) \approx \sum_{\mathbf{h}} P(\mathbf{h} | \mathbf{v}) P(v_j | \mathbf{h}) du_j$$

Now we make further approximations to deal with the exponential sum over  $\mathbf{h}$ .

**MCMC.** Given the sampling from  $P(\mathbf{h}|\mathbf{v})$  described in Section 2.2, we obtain

$$P(v_j|\mathbf{v}) \approx \frac{1}{n} \sum_{s=1}^n P(v_j|\mathbf{h}^{(s)}) du_j$$

where  $n$  is the sample size, and  $P(v_j|\mathbf{h}^{(s)})$  is computed using Eq. (6).

**Variational method.** The idea is similar to mean-field described in Section 2.2. In particular, we estimate  $\hat{h}_k = P(h_k = 1|\mathbf{v})$  using either MCMC sampling or mean-field update. The predictive distribution is approximated as

$$P(v_j|\mathbf{v}) \approx \int_{u_i \in \Omega(v_j)} P(u_i | \hat{h}_1, \hat{h}_2, \dots, \hat{h}_K)$$

where  $P(u_i | \hat{h}_1, \hat{h}_2, \dots, \hat{h}_K) = \mathcal{N}\left(u_i; \sigma_i^2 \left(\alpha_i + \sum_{k=1}^K w_{ik} \hat{h}_k\right), \sigma_i\right)$ . The computation is identical to that of Eq. (6) if we replace  $h_k$  (binary) by  $\hat{h}_k$  (real-valued).

## 2.4 Stochastic Gradient Learning with Persistent Markov Chains

Learning is based on maximising the data log-likelihood

$$\begin{aligned} \mathcal{L} &= \log P(\mathbf{v}) = \log \sum_{\mathbf{h}} \int_{\mathbf{u}} P(\mathbf{v}, \mathbf{u}, \mathbf{h}) d\mathbf{u} \\ &= \log Z(\mathbf{v}) - \log Z \end{aligned}$$

where  $P(\mathbf{v}, \mathbf{u}, \mathbf{h})$  is defined in Eq. (2) and  $Z(\mathbf{v}) = \sum_{\mathbf{h}} \int_{\mathbf{u} \in \Omega(\mathbf{v})} \Psi(\mathbf{u}, \mathbf{h}) d\mathbf{u}$ . Note that  $Z(\mathbf{v})$  includes  $Z$  as a special case when the domain  $\Omega(\mathbf{v})$  is the whole real space  $\mathbb{R}^N$ .

Recall that the model belongs to the exponential family in that we can rewrite the potential function as

$$\Psi(\mathbf{u}, \mathbf{h}) = \exp \left\{ \sum_a W_a f_a(\mathbf{u}, \mathbf{h}) \right\}$$

where  $f_a(\mathbf{u}, \mathbf{h}) \in \{u_i, u_i h_k, h_k\}_{(i,k)=(1,1)}^{(N,K)}$  is a sufficient statistic, and  $W_a \in \{\alpha_i, \gamma_k, w_{ik}\}_{(i,k)=(1,1)}^{(N,K)}$  is its associated parameter. Now the gradient of the log-likelihood has the standard form of difference of expected sufficient statistics (ESS)

$$\partial_{W_a} \mathcal{L} = \langle f_a \rangle_{P(\mathbf{u}, \mathbf{h}|\mathbf{v})} - \langle f_a \rangle_{P(\mathbf{u}, \mathbf{h})}$$

where  $P(\mathbf{u}, \mathbf{h} | \mathbf{v})$  is a truncated Gaussian RBM and  $P(\mathbf{u}, \mathbf{h})$  is the standard Gaussian RBM.

Put in common RBM-terms, there are two learning phases: the *clamped phase* in which we estimate the ESS w.r.t. the empirical distribution  $P(\mathbf{u}, \mathbf{h} | \mathbf{v})$ , and the *free phase* in which we compute the ESS w.r.t. model distribution  $P(\mathbf{u}, \mathbf{h})$ .

### 2.4.1 Persistent Markov Chains

The literature offers efficient stochastic gradient procedures to learn parameters, in which the method of [42] and its variants – the Contrastive Divergence of [13] and its persistent version of [37] – are highly effective in large-scale settings. The strategy is to update parameters after short Markov chains. Typically only the free phase requires the MCMC approximation. In our setting, on the other hand, both the clamped phase and the free phase require approximation.

Since it is possible to integrate over utilities when the binary factors are known, it is tempting to sample only the binary factors in the Rao-Blackwellisation fashion. However, here we take the advantage of the bipartite structure of the underlying RBM: the layer-wise sampling is efficient and much simpler. Once the hidden factor samples are obtained, we integrate over utilities for better numerical stability. The ESSes are the averaged over all factor samples.

For the clamped phase, we maintain one Markov chain per data instance. For memory efficiency, only the binary factor samples are stored between update steps. For the free phase, there are two strategies:

- *Contrastive chains*: one short chain is needed per data instance, but initialised from the clamped chain. That is, we discard those chains after each update.
- *Persistent chains*: free-phase chains are maintained during the course of learning, independent of the clamp-phase chains. If every data instance has the same dimensions (which they do not, in the case of missing data), we need to maintain a moderate number of chains (e.g., 20 – 100). Otherwise, we need one chain per data instance.

At each step, we collect a small number of samples and estimate the approximate distributions  $\tilde{P}(\mathbf{u}, \mathbf{h} \mid \mathbf{v})$  and  $\tilde{P}(\mathbf{u}, \mathbf{h})$ . The parameters are updated according to the stochastic gradient ascent rule

$$W_s \leftarrow W_s + \nu \left( \langle f_a \rangle_{\tilde{P}(\mathbf{u}, \mathbf{h} \mid \mathbf{v})} - \langle f_a \rangle_{\tilde{P}(\mathbf{u}, \mathbf{h})} \right)$$

where  $\nu \in (0, 1)$  is the learning rate.

### 2.4.2 Learning Thresholds

Thresholds appear only in the computation of  $Z(\mathbf{v})$  as they define the utility domain  $\Omega(\mathbf{v})$ . Let  $\bar{\Omega}(\mathbf{v})^+$  be the upper boundary of  $\Omega(\mathbf{v})$ , and  $\bar{\Omega}(\mathbf{v})^-$  the lower boundary. The gradient of the log-likelihood w.r.t. boundaries reads

$$\begin{aligned} \partial_{\bar{\Omega}(\mathbf{v})^+} \mathcal{L} &= \frac{1}{Z(\mathbf{v})} \sum_{\mathbf{h}} \partial_{\bar{\Omega}(\mathbf{v})^+} \int_{\mathbf{u} \in \Omega(\mathbf{v})} \Psi(\mathbf{h}, \mathbf{u}) d\mathbf{u} = \sum_{\mathbf{h}} P(\mathbf{u} = \bar{\Omega}(\mathbf{v})^+, \mathbf{h} \mid \mathbf{v}) \\ \partial_{\bar{\Omega}(\mathbf{v})^-} \mathcal{L} &= - \sum_{\mathbf{h}} P(\mathbf{u} = \bar{\Omega}(\mathbf{v})^-, \mathbf{h} \mid \mathbf{v}) \end{aligned}$$

Recall from Section 2.1 that the boundaries  $\bar{\Omega}(v_i = l)^-$  and  $\bar{\Omega}(v_i = l)^+$  are the lower-threshold  $\theta_{i(l-1)}$  and the upper-threshold  $\theta_{il}$ , respectively, where  $\theta_{il} = \theta_{i(l-1)} + e^{\tau_{il}} = \tau_{i1} + \sum_{m=2}^l e^{\tau_{im}}$ . Using the chain rule, we would derive the derivatives w.r.t. to  $\{\tau_{im}\}_{m=2}^{L-1}$ .

## 2.5 Handling Heterogeneous Data

We now consider the case where ordinal variables do not share the same ordinal scales, that is, we have a separate ordered set  $S_i = \{c_{i1} \prec c_{i2} \prec \dots \prec c_{iL_i}\}$  for each variable  $i$ . This requires

only slight change from the homogeneous case, e.g., by learning separate set of thresholds for each variable.

### 3 CRBM for Matrix Data

Often the data has the matrix form, i.e., a list of column vectors and we often assume columns as independent. However, this assumption is too strong in many applications. For example, in collaborative filtering where each user plays the role of a column, and each item the role of a row, a user's choice can be influenced by other users' choices (e.g., due to the popularity of a particular item), then columns are correlated. Second, it is also natural to switch the roles of the users and items and this clearly destroys the i.i.d assumption over the columns.

Thus, it is more precise to assume that an observation is *jointly* generated by both the row-wise and column-wise processes [40]. In particular, let  $d$  be the index of the data instance, each observation  $v_{di}$  is generated from an utility  $u_{di}$ . Each data instance (column)  $d$  is represented by a vector of binary hidden factors  $\mathbf{h}_d \in \{0, 1\}^K$  and each item (row)  $i$  is represented by a vector of binary hidden factors  $\mathbf{g}_i \in \{0, 1\}^S$ . Since our data matrix is usually incomplete, let us denote by  $W \in \{0, 1\}^{D \times N}$  the incidence matrix where  $W_{di} = 1$  if the cell  $(d, i)$  is observed, and  $W_{di} = 0$  otherwise. There is a single model for the whole incomplete data matrix. Every observed entry  $(d, i)$  is connected with two sets of hidden factors  $\mathbf{h}_d$  and  $\mathbf{g}_i$ . Consequently, there are  $DK + NS$  binary factor units in the entire model.

Let  $\mathbf{H} = (\{u_{di}\}_{W_{di}=1}, \{\mathbf{h}_d\}_{d=1}^D, \{\mathbf{g}_i\}_{i=1}^N)$  denote all latent variables and  $\mathbf{V} = \{v_{di}\}_{W_{di}=1}$  all visible ordinal variables. The matrix-variate model distribution has the usual form

$$P(\mathbf{V}, \mathbf{H}) = \frac{1}{Z^*} \Psi^*(\mathbf{H}) \prod_{d,i|W_{di}=1} P(v_{di} | u_{di})$$

where  $Z^*$  is the normalising constant and  $\Psi^*(\mathbf{H})$  is the product of all local potentials. More specifically,

$$\Psi^*(\mathbf{H}) = \prod_{d,i|W_{di}=1} \left( \phi_{di}(u_{di}) \prod_k \psi_{ik}(u_{di}, h_{dk}) \prod_s \varphi_{is}(u_{di}, g_{is}) \right) \left[ \prod_{d,k} \phi_k(h_{dk}) \right] \left[ \prod_{i,s} \phi_s(g_{is}) \right]$$

where  $\psi_{ik}(u_{di}, h_{dk}), \phi_k(h_{dk})$  are the same as those defined in Eq. (3), respectively, and

$$\phi_{di}(u_{di}) = \exp \left\{ -\frac{u_{di}^2}{2\sigma_{di}^2} + (\alpha_i + \beta_d)u_i \right\}; \quad \varphi_{ds}(u_{di}, g_s) = \exp \{ \omega_{ds} u_{di} g_s \}; \quad \phi_s(g_{is}) = \exp \{ \xi_s g_{is} \}$$

The ordinal model  $P(v_{di} | u_{di})$  is similar to that defined in Eq. (1) except for the thresholds, which are now functions of both the data instance and the item, that is  $\theta_{di1} = \tau_{i1} + \kappa_{d1}$  and  $\theta_{dil} = \theta_{di(l-1)} + e^{\tau_{il} + \kappa_{dl}}$  for  $l = 2, 3, \dots, L-1$ .

#### 3.1 Model Properties

It is easy to see that conditioned on the utilities, the posteriors of the binary factors are still factorisable. Likewise, given the factors, the utilities are univariate Gaussian



$$\begin{aligned}
P(u_{di} \mid \mathbf{h}_d, \mathbf{g}_i) &= \mathcal{N}(\mu_{di}^*(\mathbf{h}_d, \mathbf{g}_i), \sigma_{di}^2) \\
P(u_{di} \mid \mathbf{h}_d, \mathbf{g}_i, v_{di}) &\propto \mathbb{I}[u_{di} \in \Omega(v_{di})] P(u_{di} \mid \mathbf{h}_d, \mathbf{g}_i)
\end{aligned}$$

where  $\Omega(v_{di})$  is the domain defined by the thresholds at the level  $l = v_{di}$ , and the mean structure is

$$\mu_{di}^*(\mathbf{h}_d, \mathbf{g}_i) = \sigma_{di}^2 \left( \alpha_i + \beta_d + \sum_{k=1}^K w_{ik} h_{dk} + \sum_{s=1}^S \omega_{ds} g_{is} \right) \quad (12)$$

Previous inference tricks can be re-used by noting that for each column (i.e., data instance), we still enjoy the Gaussian RBM when conditioned on other columns. The same holds for rows (i.e., items).

### 3.2 Stochastic Learning with Structured Mean-Fields

Although it is possible to explore the space of the whole model using Gibbs sampling and use the short MCMC chains as before, here we resort to structured mean-field methods to exploit the modularity in the model structure. The general idea is to alternate between the column-wise and the row-wise conditional processes:

- In the *column-wise* process, we estimate item-specific factor posteriors  $\{\hat{\mathbf{g}}_i\}_{i=1}^N$ , where  $\hat{g}_{is} \leftarrow P(g_{is} = 1 \mid (v_{di})_{di \mid W_{id}=1})$  and use them *as if* the item-specific factors  $(\mathbf{g}_i)_{i=1}^N$  are given. For example, the mean structure in Eq. (12) now has the following form

$$\mu_{di}^*(\mathbf{h}_d, \hat{\mathbf{g}}_i) = \sigma_{di}^2 \left( \alpha_i + \beta_d + \sum_{k=1}^K w_{ik} h_{dk} + \sum_{s=1}^S \omega_{ds} \hat{g}_{is} \right)$$

which is essentially the mean structure in Eq. (5) when  $\beta_d + \sum_{s=1}^S \omega_{ds} \hat{g}_{is}$  is absorbed into  $\alpha_i$ . Conditioned on the estimated posteriors, the data likelihood is now factorisable  $\prod_d P(\mathbf{v}_{d\bullet} \mid \{\hat{\mathbf{g}}_i\}_{i=1}^N)$ , where  $\mathbf{v}_{d\bullet}$  denotes the observations of the  $d$ -th data instance.

- Similarly, in the *row-wise* process we estimate data-specific posteriors  $\{\hat{\mathbf{h}}_d\}_{d=1}^D$ , where  $\hat{h}_{dk} = P(h_{dk} = 1 \mid (v_{di})_{W_{id}=1})$  and use them *as if* the data-specific factors  $(\mathbf{h}_d)_{d=1}^D$  are given. The data likelihood has the form  $\prod_i P(\mathbf{v}_{\bullet i} \mid \{\hat{\mathbf{h}}_d\}_{d=1}^D)$ , where  $\mathbf{v}_{\bullet i}$  denotes the observations of the  $i$ -th item.

At each step, we then improve the conditional data likelihood using the gradient technique described in Section 2.4, e.g., by running through the whole data once.

#### 3.2.1 Online Estimation of Posteriors

The structured mean-fields technique requires the estimation of the factor posteriors. To reduce computation, we propose to treat the trajectory of the factor posteriors during learning as a *stochastic process*. This suggests a simple smoothing method, e.g., at step  $t$ :

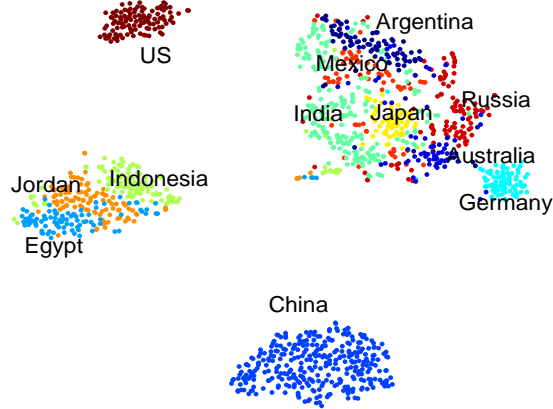


Figure 2: Visualisation of world’s opinions in 2008 by projecting latent posteriors  $\hat{\mathbf{h}} = (P(h_1^1|\mathbf{v}), P(h_2^1|\mathbf{v}), \dots, P(h_K^1|\mathbf{v}))$  on 2D using t-SNE [41], where  $h_k^1$  is a shorthand for  $h_k = 1$ . Best viewed in colours. .

$$\hat{\mathbf{h}}_d^{(t)} \leftarrow \eta \hat{\mathbf{h}}_d^{(t-1)} + (1 - \eta) P(h_{dk} = 1 | \mathbf{u}_d^{(t)})$$

where  $\eta \in (0, 1)$  is the *smoothing factor*, and  $\mathbf{u}_d^{(t)}$  is a utility sample in the *clamped phase*. This effectively imposes an exponential decay to previous samples. The estimation of  $\eta$  would be of interest in its own right, but we would empirically set  $\eta \in (0.5, 0.9)$  and do not pursue the issue further.

## 4 Experiments

In this section, we demonstrate how CRBM can be useful in real-world data analysis tasks. To monitor learning progress, we estimate the data pseudo-likelihood  $P(v_i | \mathbf{v}_{-i})$ . For simplicity, we treat  $v_i$  as if it is not in  $\mathbf{v}$  and replace  $\mathbf{v}_i$  by  $\mathbf{v}$ . This enables us to use the same predictive methods in Section 2.3. See Fig. 3(a) for an example of the learning curves. To sample from the truncated Gaussian, we employ methods described in [30], which is more efficient than standard rejection sampling techniques. Mapping parameters  $\{w_{ik}\}$  are initialised randomly, bias paramters are from zeros, and thresholds  $\{\theta_{il}\}$  are spaced evenly at the begining.

### 4.1 Global Attitude Analysis: Latent Profile Discovery

In this experiments we validate the capacity to discover meaningful latent profiles from people’s opinions about their life and the social/political conditions in their country and around the world. We use the public world-wide survey by PewResearch Centre<sup>3</sup> in 2008 which interviewed 24, 717 people from 24 countries. After re-processing, we keep 165 ordinal responses per respondent.

<sup>3</sup><http://pewresearch.org/>

Example questions are: “(Q1) [...] how would you describe your day today—has it been a typical day, a particularly good day, or a particularly bad day?”, “(Q5) [...] over the next 12 months do you expect the economic situation in our country to improve a lot, improve a little, remain the same, worsen a little or worsen a lot?”.

The data is heterogeneous since question types are different (see Section 2.5). For this we use a vector-based CRBM with  $K = 50$  hidden units. After model fitting, we obtain a posterior vector  $\hat{\mathbf{h}} = (P(h_1^1|\mathbf{v}), P(h_2^1|\mathbf{v}), \dots, P(h_K^1|\mathbf{v}))$ , which is then used as the representation of the respondent’s latent profile. For visualisation, we project this vector onto the 2D plane using a locality-preserving dimensionality reduction method known as *t-SNE*<sup>4</sup> [41]. The opinions of citizens of 12 countries are depicted in Fig. 2. This clearly reveals how cultures (e.g., Islamic and Chinese) and nations (e.g., the US, China, Latin America) see the world.

## 4.2 Collaborative Filtering: Matrix Completion

We verify our models on three public rating datasets: MovieLens<sup>5</sup> – containing 1 million ratings by 6 thousand users on nearly 4 thousand movies; Dating<sup>6</sup> – consisting of 17 million ratings by 135 thousand users on nearly 169 thousand profiles; and Netflix<sup>7</sup> – 100 millions ratings by 480 thousand users on nearly 18 thousand movies. The Dating ratings are on the 10-point scale and the other two are on the 5-star scale. We then transform the Dating ratings to the 5-point scale for uniformity. For each data we remove those users with less than 30 ratings, 5 of which are used for tuning and stopping criterion, 10 for testing and the rest for training. For MovieLens and Netflix, we ensure that rating timestamps are ordered from training, to validation to testing. For the Dating dataset, the selection is at random.

For comparison, we implement state-of-the-art methods in the field, including: Matrix Factorisation (MF) with Gaussian assumption [34], MF with cumulative ordinal assumption [16] (without item-item neighbourhood), and RBM with multinomial assumption [35]. For prediction in the CRBM, we employ the variational method (Section 11). The training and testing protocols are the same for all methods: Training stops where there is no improvement on the likelihood of the validation data. Two popular performance metrics are reported on the test data: the *root-mean square error* (RMSE), the *mean absolute error* (MAE). Prediction for ordinal MF and RBMs is a numerical mean in the case of RMSE:  $v_j^{RMSE} = \sum_{l=1}^L P(v_j = l|\mathbf{v})l$ , and an MAP estimation in the case of MAE:  $v_j^{MAE} = \arg \max_l P(v_j = l|\mathbf{v})$ .

Fig. 3(a) depicts the learning curve of the vector-based and matrix-based CRBMs, and Fig. 3(b) shows their predictive performance on test datasets. Clearly, the effect of matrix treatment is significant. Tables 1,2,3 report the performances of all methods on the three datasets. The (matrix) CRBM are often comparable with the best rivals on the RMSE scores and are competitive against all others on the MAE.

## 5 Related Work

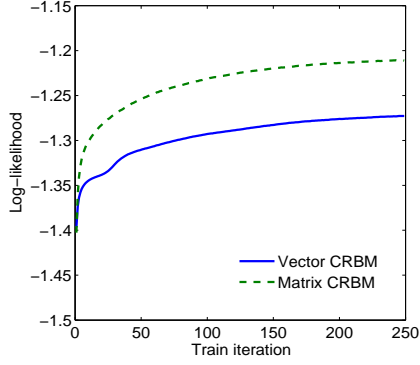
This work partly belongs to the thread of research that extends RBMs for a variety of data types, including *categories* [35], *counts* [10, 33, 32], *bounded* variables [19] and a mixture of these types

<sup>4</sup>Note that the t-SNE does not do clustering, it tries only to map from the input to the 2D so that local properties of the data is preserved.

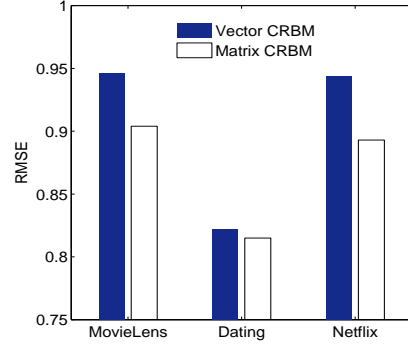
<sup>5</sup><http://www.grouplens.org/node/12>

<sup>6</sup><http://www.occamlab.com/petricek/data/>

<sup>7</sup><http://netflixprize.com/>



(a) Monitoring pseudo-likelihood in training



(b) RMSE on test data

Figure 3: Vector versus matrix CRBMs, where  $K = 50$ .

	$K = 50$		$K = 100$		$K = 200$	
	<i>RMSE</i>	<i>MAE</i>	<i>RMSE</i>	<i>MAE</i>	<i>RMSE</i>	<i>MAE</i>
Gaussian Matrix Fac.	0.914	0.720	0.911	0.719	0.908	0.716
Ordinal Matrix Fac.	<b>0.904</b>	0.682	<b>0.902</b>	0.682	<b>0.902</b>	0.680
Multinomial RBM	0.928	0.711	0.926	0.707	0.928	0.708
Matrix Cumul. RBM	<b>0.904</b>	<b>0.666</b>	0.904	<b>0.662</b>	0.906	<b>0.664</b>

Table 1: Results on MovieLens (the smaller the better).

	$K = 50$		$K = 100$		$K = 200$	
	<i>RMSE</i>	<i>MAE</i>	<i>RMSE</i>	<i>MAE</i>	<i>RMSE</i>	<i>MAE</i>
Gaussian Matrix Fac.	0.852	0.596	0.848	0.592	0.840	0.586
Ordinal Matrix Fac.	0.857	0.511	0.854	0.507	0.849	0.502
Multinomial RBM	<b>0.815</b>	0.483	<b>0.794</b>	0.470	0.787	0.463
Matrix Cumul. RBM	<b>0.815</b>	<b>0.475</b>	0.799	<b>0.461</b>	0.794	<b>0.458</b>

Table 2: Results on Dating (the smaller the better).

	$K = 50$		$K = 100$	
	<i>RMSE</i>	<i>MAE</i>	<i>RMSE</i>	<i>MAE</i>
Gaussian Matrix Fac.	<b>0.890</b>	0.689	0.888	0.688
Ordinal Matrix Fac.	0.904	0.658	0.902	0.657
Multinomial RBM	0.894	0.659	<b>0.887</b>	0.650
Matrix Cumul. RBM	0.893	<b>0.641</b>	0.892	<b>0.640</b>

Table 3: Results on Netflix (the smaller the better).

[38]. Gaussian RBMs have been only used for continuous variables [13, 25] – thus our use for ordinal variables is novel. There has also been recent work extending Gaussian RBMs to better model highly correlated input variables [28, 8]. For ordinal data, to the best of our knowledge, the first RBM-based work is [40], which also contains a treatment of matrix-wise data. However, their work indeed models multinomial data with knowledge of orders rather than modelling the ordinal nature directly. The result is that it is over-parameterised but less efficient and does not offer any underlying generative mechanism for ordinal data.

Ordinal data has been well investigated in statistical sciences, especially quantitative social studies, often under the name of *ordinal regression*, which refers to single ordinal output given a set of input covariates. The most popular method is by [24] which examines the level-wise cumulative distributions. Another well-known treatment is the sequential approach, also known as continuation ratio [23], in which the ordinal generation process is considered stepwise, starting from the lowest level until the best level is chosen. For reviews of recent development, we refer to [22]. In machine learning, this has attracted a moderate attention in the past decade [12, 6, 7, 3], adding machine learning flavours (e.g., large-margins) to existing statistical methods.

Multivariate ordinal variables have also been studied for several decades [1]. The most common theme is the assumption of the latent multivariate normal distribution that generates the ordinal observations, often referred to as *multivariate probit* models [5, 11, 17, 27, 15, 4]. The main problem with this setting is that it is only feasible for problems with small dimensions. Our treatment using RBMs offer a solution for large-scale settings by transferring the low-order interactions among the Gaussian variables onto higher-order interactions through the hidden binary layer. Not only this offers much faster inference, it also enables automatic discovery of latent aspects in the data.

For matrix data, the most well-known method is perhaps matrix factorisation [21, 29, 34]. However, this method assumes that the data is normally distributed, which does not meet the ordinal characteristics well. Recent research has attempted to address this issue [26, 16, 39]. In particular, [26, 16] adapt cumulative models of [24], and [39] tailors the sequential models of [23] for task.

## 6 Conclusion

We have presented CRBM, a novel probabilistic model to handle vector-variate and matrix-variate ordinal data. The model is based on Gaussian restricted Boltzmann machines and we present the model architecture, learning and inference procedures. We show that the model is useful in profiling opinions of people across cultures and nations. The model is also competitive against state-of-art methods in collaborative filtering using large-scale public datasets. Thus our work enriches the RBMs, and extends their use on multivariate ordinal data in diverse applications.

## References

- [1] J.A. Anderson and JD Pemberton. The grouped continuous model for multivariate ordered categorical variables and covariate adjustment. *Biometrics*, pages 875–885, 1985.
- [2] Bo Chen Benjamin Marlin, Kevin Swersky and Nando de Freitas. Inductive Principles for Restricted Boltzmann Machine Learning. In *Proceedings of the 13rd International Con-*

*ference on Artificial Intelligence and Statistics*, Chia Laguna Resort, Sardinia, Italy, May 2010.

- [3] J.S. Cardoso and J.F.P. da Costa. Learning to classify ordinal data: the data replication method. *Journal of Machine Learning Research*, 8(1393-1429):6, 2007.
- [4] P. Chagneau, F. Mortier, N. Picard, and J.N. Bacro. A hierarchical bayesian model for spatial prediction of multivariate non-gaussian random fields. *Biometrics*, 2010.
- [5] S. Chib and E. Greenberg. Analysis of multivariate probit models. *Biometrika*, 85(2):347–361, 1998.
- [6] W. Chu and Z. Ghahramani. Gaussian processes for ordinal regression. *Journal of Machine Learning Research*, 6(1):1019, 2006.
- [7] W. Chu and S.S. Keerthi. Support vector ordinal regression. *Neural computation*, 19(3):792–815, 2007.
- [8] A. Courville, J. Bergstra, and Y. Bengio. A spike and slab restricted Boltzmann machine. In *Proceedings of the 14th International Conference on Artificial Intelligence and Statistics (AISTATS)*, volume 15. Fort Lauderdale, USA, 2011.
- [9] Y. Freund and D. Haussler. Unsupervised learning of distributions on binary vectors using two layer networks. *Advances in Neural Information Processing Systems*, pages 912–919, 1993.
- [10] P.V. Gehler, A.D. Holub, and M. Welling. The rate adapting Poisson model for information retrieval and object recognition. In *Proceedings of the 23rd international conference on Machine learning*, pages 337–344. ACM New York, NY, USA, 2006.
- [11] Leonardo Grilli and Carla Rampichini. Alternative specifications of multivariate multilevel probit ordinal response models. *Journal of Educational and Behavioral Statistics*, 28(1):31–44, 2003.
- [12] R. Herbrich, T. Graepel, and K. Obermayer. Large margin rank boundaries for ordinal regression. *Advances in Neural Information Processing Systems*, pages 115–132, 1999.
- [13] G.E. Hinton. Training products of experts by minimizing contrastive divergence. *Neural Computation*, 14:1771–1800, 2002.
- [14] G.E. Hinton and R.R. Salakhutdinov. Reducing the dimensionality of data with neural networks. *Science*, 313(5786):504–507, 2006.
- [15] I. Jeliazkov, J. Graves, and M. Kutzbach. Fitting and comparison of models for multivariate ordinal outcomes. *Advances in Econometrics*, 23:115–156, 2008.
- [16] Y. Koren and J. Sill. OrdRec: an ordinal model for predicting personalized item rating distributions. In *Proceedings of the fifth ACM conference on Recommender systems*, pages 117–124. ACM, 2011.
- [17] A. Kottas, P. Müller, and F. Quintana. Nonparametric Bayesian modeling for multivariate ordinal data. *Journal of Computational and Graphical Statistics*, 14(3):610–625, 2005.

- [18] H. Larochelle and Y. Bengio. Classification using discriminative restricted Boltzmann machines. In *Proceedings of the 25th international conference on Machine learning*, pages 536–543. ACM, 2008.
- [19] Q.V. Le, W.Y. Zou, S.Y. Yeung, and A.Y. Ng. Learning hierarchical invariant spatio-temporal features for action recognition with independent subspace analysis. *CVPR*, 2011.
- [20] N. Le Roux and Y. Bengio. Representational power of restricted Boltzmann machines and deep belief networks. *Neural Computation*, 20(6):1631–1649, 2008.
- [21] D.D. Lee and H.S. Seung. Learning the parts of objects by non-negative matrix factorization. *Nature*, 401(6755):788–791, 1999.
- [22] I. Liu and A. Agresti. The analysis of ordered categorical data: an overview and a survey of recent developments. *TEST*, 14(1):1–73, 2005.
- [23] R.D. Mare. Social background and school continuation decisions. *Journal of the American Statistical Association*, pages 295–305, 1980.
- [24] P. McCullagh. Regression models for ordinal data. *Journal of the Royal Statistical Society. Series B (Methodological)*, pages 109–142, 1980.
- [25] A.R. Mohamed and G. Hinton. Phone recognition using restricted Boltzmann machines. In *Acoustics Speech and Signal Processing (ICASSP), 2010 IEEE International Conference on*, pages 4354–4357. IEEE, 2010.
- [26] U. Paquet, B. Thomson, and O. Winther. A hierarchical model for ordinal matrix factorization. *Statistics and Computing*, pages 1–13, 2011.
- [27] J. Podani. Multivariate exploratory analysis of ordinal data in ecology: pitfalls, problems and solutions. *Journal of Vegetation Science*, 16(5):497–510, 2005.
- [28] M.A. Ranzato and G.E. Hinton. Modeling pixel means and covariances using factorized third-order Boltzmann machines. In *CVPR*, pages 2551–2558. IEEE, 2010.
- [29] J.D.M. Rennie and N. Srebro. Fast maximum margin matrix factorization for collaborative prediction. In *Proceedings of the 22nd International Conference on Machine Learning (ICML)*, pages 713–719, Bonn, Germany, 2005.
- [30] C.P. Robert. Simulation of truncated normal variables. *Statistics and computing*, 5(2):121–125, 1995.
- [31] R. Salakhutdinov and G. Hinton. Deep Boltzmann Machines. In *Proceedings of The Twelfth International Conference on Artificial Intelligence and Statistics (AISTATS’09, volume 5*, pages 448–455, 2009.
- [32] R. Salakhutdinov and G. Hinton. Replicated softmax: an undirected topic model. *Advances in Neural Information Processing Systems*, 22, 2009.
- [33] R. Salakhutdinov and G. Hinton. Semantic hashing. *International Journal of Approximate Reasoning*, 50(7):969–978, 2009.

- [34] R. Salakhutdinov and A. Mnih. Probabilistic matrix factorization. *Advances in neural information processing systems*, 20:1257–1264, 2008.
- [35] R. Salakhutdinov, A. Mnih, and G. Hinton. Restricted Boltzmann machines for collaborative filtering. In *Proceedings of the 24th International Conference on Machine Learning (ICML)*, pages 791–798, 2007.
- [36] P. Smolensky. Information processing in dynamical systems: Foundations of harmony theory. *Parallel distributed processing: Explorations in the microstructure of cognition*, 1:194–281, 1986.
- [37] T. Tieleman. Training restricted Boltzmann machines using approximations to the likelihood gradient. In *Proceedings of the 25th international conference on Machine learning*, pages 1064–1071. ACM, 2008.
- [38] T. Tran, D.Q. Phung, and S. Venkatesh. Mixed-variate restricted Boltzmann machines. In *Proc. of 3rd Asian Conference on Machine Learning (ACML)*, Taoyuan, Taiwan, 2011.
- [39] T. Tran, D.Q. Phung, and S. Venkatesh. Sequential decision approach to ordinal preferences in recommender systems. In *Proc. of the 26th AAAI Conference*, Toronto, Ontario, Canada, 2012.
- [40] T.T. Truyen, D.Q. Phung, and S. Venkatesh. Ordinal Boltzmann machines for collaborative filtering. In *Twenty-Fifth Conference on Uncertainty in Artificial Intelligence (UAI)*, Montreal, Canada, June 2009.
- [41] L. van der Maaten and G. Hinton. Visualizing data using t-SNE. *Journal of Machine Learning Research*, 9(2579-2605):85, 2008.
- [42] L. Younes. Parametric inference for imperfectly observed Gibbsian fields. *Probability Theory and Related Fields*, 82(4):625–645, 1989.